

Joint Cache-Channel Coding over Erasure Broadcast Channels

Roy Timo and Michèle Wigger

Abstract—We consider a cache-aided communications system in which a transmitter communicates with many receivers over an erasure broadcast channel. The system serves as a basic model for communicating on-demand content during periods of high network congestion, where some content can be pre-placed in local caches near the receivers. We formulate the cache-aided communications problem as a joint cache-channel coding problem, and characterise some information-theoretic tradeoffs between reliable communications rates and cache sizes. We show that if the receivers experience different channel qualities, then using unequal cache sizes and joint cache-channel coding improves system efficiency.

I. INTRODUCTION

Consider a network with one transmitter and many receivers. Imagine that the transmitter has a library of *messages* (or, *data files*), and suppose that each receiver will request and download a message during a period of high network congestion. In such settings, it is advantageous to move traffic away from the congested period using *caching*. The basic idea of caching is that the transmitter sends and stores “parts” of the library in local *cache memories* near the receivers beforehand, during periods with low network traffic. The caches provide this data directly to the receivers, so that less data needs to be sent during the congested period.

The above problem is relevant to video-streaming services, where content providers pre-place data in clients’ caches (or, on servers near the clients), with the goal of improving latency and rate performance in high demand periods. The content provider typically does not know in advance which specific movies the clients will request, and thus the cached data cannot depend on the clients’ specific demands.

Let us call the pre-placement of data in caches the *caching phase*, and the remaining communications phase the *delivery phase*. Cache memories are typically much smaller than the library, and the caching phase occurs before the receivers’ demands are known. A key engineering challenge is, therefore, to carefully choose and cache only that data which is most useful during the delivery phase. That is, one should cache data that minimises the rate needed to complete the delivery-phase downloads for any feasible receiver demands.

Cache-aided communications systems have received significant attention in the information-theoretic literature in recent years, and those works most closely related to this paper are [1]–[11]. With the exception of [11], these works assume that the delivery phase takes place over a single rate-limited multicast

noiseless channel (a bit-pipe) that connects the transmitter to every receiver. In practice, however, the communications medium is sometimes better modelled by a noisy broadcast channel (BC). This scenario is considered in [11], where the BC is essentially a set of parallel links with different qualities to the various receivers, which models a wireless fading BC.

This paper takes a similar approach to that of [11], and we assume that the delivery phase takes place over a memoryless erasure BC. However, in contrast to [11], we assume that the caching phase takes place over error-free pipes. The motivation for this simplified assumption is that the caching phase typically occurs during periods of low network-congestion, where network resources are not a limiting factor.

Our main contribution in this paper is a joint cache-channel coding scheme for the described setup for general demands, and a characterization of the *capacity-memory* region when the receivers wish to learn the same message. Our results show that when the receivers experience different erasure probabilities (different channel qualities), then

- it is beneficial to employ unequal cache sizes at the receivers (larger cache memories at weaker receivers, and smaller cache memories at strong receivers); and
- joint cache-channel coding techniques can provide significant gains over separated cache and channel coding.

Allocating larger cache memories to the weaker receivers is quite natural because one then needs to communicate less data over noisier channels (see also [11]). Interestingly, there is an additional benefit to asymmetric caches that arises when joint cache-channel coding is used during the delivery phase. The basic idea is as follows: Consider a degraded BC communications scenario (such as the erasure BC) with separate cache and channel coding. Here a stronger receiver can decode all the data that is sent to a weaker receiver during the delivery phase. In fact, the strong receiver could decode even more data, but it is limited by the weaker receiver. Now suppose that part of the message intended for the stronger receiver is stored within the weaker receiver’s cache: one can freely piggyback this part of the stronger receiver’s message on the message intended for the weaker receiver. The weaker receiver is not penalised because it knows what data is being piggybacked on its desired message, and its channel decoder can still resolve its desired message. While, simultaneously, the stronger receiver has decoded something about its desired message and therefore we have improved efficiently. Thus, thanks to the weaker receiver’s cache and a simple joint cache-channel coding scheme, we can send additional data to stronger receivers without any extra cost, i.e.,

R. Timo is with TU München (roy.timo@tum.de), and M. Wigger is with Telecom ParisTech (michele.wigger@telecom-paristech.fr). This work was supported by the Alexander von Humboldt Foundation.

extra rate-constraints. This additional benefit of asymmetric cache memories was not observed in [11], because a separate source-channel coding scheme was used for the delivery phase.

II. PROBLEM DEFINITION

A. Message library and feasible receiver demands

We have a transmitter, K receivers and a library with D messages W_1, \dots, W_D . The d -th message in the library W_d is independent of all other messages and uniform on¹

$$\{0, 1, \dots, 2^{nR_d} - 1\},$$

where $R_d \geq 0$ is its rate and n is the transmission blocklength. We represent a particular combination of receivers' demands by a tuple $\mathbf{d} = (d_1, \dots, d_K) \in \{1, \dots, D\}^K$. That is, \mathbf{d} represents the situation where receiver 1 demands (i.e., requests and downloads) message W_{d_1} , receiver 2 demands W_{d_2} , and so on. Let

$$\mathcal{D} \subseteq \{1, \dots, D\}^K.$$

denote the *feasible set* of all possible receiver demands. The feasible set \mathcal{D} is known to the transmitter and receivers during the caching and delivery phases, but the specific demand tuple \mathbf{d} chosen from \mathcal{D} is only revealed for the delivery phase.

B. Caching phase

For each receiver $k \in \{1, \dots, K\}$, the size of its cache is described by a nonnegative integer \mathcal{M}_k , see (1) below. The transmitter sends

$$\mathbb{Z}_k := g_k(W_1, \dots, W_D),$$

to receiver k 's cache, where $g_k : \prod_{d=1}^D \{0, 1, \dots, 2^{nR_d} - 1\} \rightarrow \mathcal{Z}_k$. such that

$$\log |\mathcal{Z}_k| \leq 2^{n\mathcal{M}_k}. \quad (1)$$

The caching phase occurs during a low congestion period, and we assume that \mathbb{Z}_k is reliably conveyed to receiver k 's cache (for each $k \in \{1, \dots, K\}$).

C. Erasure Broadcast Channel Model

The delivery phase occurs during a high congestion period, which we model by an *erasure BC* with input alphabet $\mathcal{X} := \{0, 1\}^F$. Here $F \geq 0$ is a fixed positive integer, and each $x \in \mathcal{X}$ is an F -bit packet. Due to congestion, some packets may be lost when, for example, router buffers overload. We denote the event of a lost packet with the *erasure symbol* Δ , and the BC's output alphabet by $\mathcal{Y} := \mathcal{X} \cup \{\Delta\}$ (the same alphabet is used for all receivers). Fix

$$1 \geq \delta_1 \geq \delta_2 \geq \dots \geq \delta_K \geq 0.$$

Let $Q(y_1, \dots, y_K | x) := \mathbb{P}[(Y_1, \dots, Y_K) = (y_1, \dots, y_K) | X = x]$ be any probability law for the memoryless BC with marginals

$$\mathbb{P}[Y_k = y_k | X = x] = \begin{cases} 1 - \delta_k & \text{if } y_k = x \\ \delta_k & \text{if } y_k = \Delta \\ 0 & \text{otherwise} \end{cases}, \quad \forall k.$$

¹To simplify notation and help elucidate our main ideas, we assume throughout the paper that 2^{nR_d} is an integer.

For our purpose only these marginal probabilities are relevant.

We discuss a caching system in the next section that is built on separate cache and channel codes, and, for this reason, it is useful to recall the degraded message set capacity region for Q . A channel-coding rate tuple $(R_{\{1, \dots, K\}}, R_{\{2, \dots, K\}}, \dots, R_{\{K\}})$ is said to be achievable on Q if the following holds: For any $\epsilon > 0$ there exists an encoder and K -decoders such that, for all k , the transmitter can send $(R_{\{k, \dots, K\}} - \epsilon)$ information bits per channel use to every receiver in the set $\{k, k+1, \dots, K\}$ with an average probability of error less than ϵ . The set of all achievable rates — the *degraded message set capacity region* \mathcal{C}^\dagger — is given by the next proposition. The proposition can be distilled from [12], and we omit these details.

Proposition 1:

$$\mathcal{C}^\dagger = \left\{ (R_{\{1, \dots, K\}}, R_{\{2, \dots, K\}}, \dots, R_{\{K\}}) : \sum_{k=1}^K \frac{R_{\{k, \dots, K\}}}{F(1 - \delta_k)} \leq 1, R_{\{k, \dots, K\}} \geq 0, \quad \forall k \right\}.$$

D. Delivery phase

For each feasible demand $\mathbf{d} \in \mathcal{D}$, let

$$f_{\mathbf{d}} : \prod_{d'=1}^D \{0, 1, \dots, 2^{nR_{d'}} - 1\} \rightarrow \mathcal{X}^n$$

denote the corresponding encoder at the transmitter. Given $\mathbf{d} \in \mathcal{D}$ and the library (W_1, \dots, W_D) , the transmitter sends

$$X^n := f_{\mathbf{d}}(W_1, \dots, W_D), \quad (2)$$

where $X^n = (X_1, \dots, X_n)$. Receiver k observes $Y_k^n = (Y_{k,1}, \dots, Y_{k,n})$ according to the memoryless law Q . Let

$$\varphi_{k,\mathbf{d}} : \mathcal{Y}^n \times \mathcal{Z}_k \rightarrow \{0, 1, \dots, 2^{nR_{d_k}} - 1\} \quad (3)$$

denote the decoder at receiver k . Given demands $\mathbf{d} \in \mathcal{D}$, cache content \mathbb{Z}_k and channel outputs Y_k^n , receiver k outputs

$$\hat{W}_k := \varphi_{k,\mathbf{d}}(Y_k^n, \mathbb{Z}_k)$$

as its reconstruction of the d_k -th message W_{d_k} .

E. Achievable rate-memory tuples

Let

$$P_e := \mathbb{P} \left[\bigcup_{\mathbf{d} \in \mathcal{D}} \bigcup_{k=1}^K \{\hat{W}_k \neq W_{d_k}\} \right]$$

denote the probability of error at any receiver for any feasible demand. We call the collection of all encoders and decoders,

$$\{g_1, g_2, \dots, g_K\} \text{ and } \{f_{\mathbf{d}}, \varphi_{1,\mathbf{d}}, \varphi_{2,\mathbf{d}}, \dots, \varphi_{K,\mathbf{d}}\}_{\mathbf{d} \in \mathcal{D}},$$

an $(n, R_1, \dots, R_D, \mathcal{M}_1, \dots, \mathcal{M}_K)$ -code.

We say that a rate-memory tuple $(R_1, \dots, R_D, \mathcal{M}_1, \dots, \mathcal{M}_K)$ is *achievable* if for any $\epsilon > 0$ there exists a sufficiently large blocklength n and an $(n, R_1, \dots, R_D, \mathcal{M}_1, \dots, \mathcal{M}_K)$ -code with $P_e \leq \epsilon$.

Definition 1: We define the *capacity-memory region* \mathcal{C} to be the closure of the set of all achievable rate-memory tuples.

The main problem of interest in this paper is to determine the capacity-memory region \mathcal{C} for a given erasure BC Q .

III. MOTIVATING EXAMPLES

We now demonstrate the potential of unequal cache memories and joint cache-channel coding with three examples. Fix $K = 2$; $\mathcal{D} = \{1, \dots, D\}^2$; $R_d = R$ for all d ; and

$$\delta_1 = 4/5 \quad \text{and} \quad \delta_2 = 1/5. \quad (4)$$

A. Coded caching with symmetric caches

Suppose that $\mathcal{M}_1 = \mathcal{M}_2 = \mathcal{M}$, and

$$\alpha := \mathcal{M}/R \in [0, D/2]. \quad (5)$$

Split each message W_d in the library into three sub-messages,

$$W_d = (W_d^{(c1)}, W_d^{(c2)}, W_d^{(u)}),$$

of rates \mathcal{M}/D , \mathcal{M}/D , and $R - 2\mathcal{M}/D$.

Caching phase: Store the sub-messages

$$(W_1^{(c1)}, \dots, W_D^{(c1)}) \quad \text{and} \quad (W_1^{(c2)}, \dots, W_D^{(c2)})$$

in the caches of receiver 1 and 2's respectively.

Delivery phase: The transmitter sends

$$W_{d_1}^{(c2)} \oplus W_{d_2}^{(c1)}, \quad (6)$$

as a common message to both receivers, where the addition is modulo $2^{n(\mathcal{M}/D)}$. It then sends $W_{d_1}^{(u)}$ as a private message to receiver 1 and $W_{d_2}^{(u)}$ as a private message to receiver 2. Notice that receiver 1 can recover W_{d_1} from the common message and $W_{d_1}^{(u)}$, while receiver 2 can recover W_{d_2} from the common message and $W_{d_2}^{(u)}$. We use a good channel code to communicate these messages over the BC.

Achievable rate-memory tuples: Proposition 1 asserts that the common message (6) and $W_{d_1}^{(u)}$ can be decoded by both receivers and $W_{d_2}^{(u)}$ can be decoded by receiver 2 whenever

$$\frac{R - \frac{\mathcal{M}}{D}}{F(1 - \delta_1)} + \frac{R - \frac{2\mathcal{M}}{D}}{F(1 - \delta_2)} \leq 1. \quad (7)$$

On substituting (4), the inequality (7) simplifies to

$$R \leq \frac{4}{5}F(1 - \delta_1) + \frac{6}{5}\frac{\mathcal{M}}{D}. \quad (8)$$

All rate-memory tuples $(R, \dots, R, \mathcal{M}, \dots, \mathcal{M})$, with R and \mathcal{M} satisfying (5) and (8), are achievable.

B. Separate cache-channel coding and asymmetric caches

Now suppose that we have asymmetric caches $\mathcal{M}_1 = 2\mathcal{M}$ and $\mathcal{M}_2 = 0$ for some \mathcal{M} satisfying (5). The total cache memory available at both receivers remains unchanged, only now the memory at receiver 2 has been reallocated to receiver 1.

Split each message W_d into two sub-messages,

$$W_d = (W_d^{(c1)}, W_d^{(u)}) \quad (9)$$

with rates $2\mathcal{M}/D$ and $R - (2\mathcal{M}/D)$ respectively.

Caching phase: Store $(W_1^{(c1)}, \dots, W_D^{(c1)})$ in receiver 1's cache.

Delivery phase: We use a good channel code for Proposition 1 to reliably communicate the above sub-messages. The transmitter sends $W_{d_1}^{(u)}$ as a common message to both receivers

(although it is only used by receiver 1), and it sends $W_{d_2}^{(c1)}$ and $W_{d_2}^{(u)}$ as a private message to receiver 2.

Achievable rate-memory tuples: Proposition 1 asserts that reliable communication is possible if

$$\frac{R - \frac{2\mathcal{M}}{D}}{F(1 - \delta_1)} + \frac{R}{F(1 - \delta_2)} \leq 1. \quad (10)$$

On substituting (4), the inequality (10) simplifies to

$$R \leq \frac{4}{5}F(1 - \delta_1) + \frac{8}{5}\frac{\mathcal{M}}{D}. \quad (11)$$

All rate-memory tuples $(R, \dots, R, \mathcal{M}, \dots, \mathcal{M})$, with R and \mathcal{M} satisfying (5) and (11) are achievable.

C. Joint cache-channel coding and asymmetric caches

As in Section III-B: Let $\mathcal{M}_2 = 0$ and $\mathcal{M}_1 = 2\mathcal{M}$, for some \mathcal{M} satisfying (5), and split each message W_d into two sub-messages (9) with rates $2\mathcal{M}/D$ and $R - (2\mathcal{M}/D)$ respectively.

Caching phase: Store $(W_1^{(c1)}, \dots, W_D^{(c1)})$ at receiver 1.

Delivery phase: Transmission takes place in two phases using timesharing. First phase of length $\beta_1 n$, for some $\beta_1 \in [0, 1]$: The transmitter sends

$$(W_{d_1}^{(u)}, W_{d_2}^{(c1)})$$

as a common message to both receivers using a joint cache-channel code. Second phase of length $(1 - \beta_1)n$: The transmitter sends $W_{d_2}^{(u)}$ to receiver 2 using a point-to-point channel code. Receiver 1 tries to decode $W_{d_1}^{(u)}$ and receiver 2 tries to decode $(W_{d_1}^{(u)}, W_{d_2}^{(c1)}, W_{d_2}^{(u)})$. A key observation here is that $W_{d_2}^{(c1)}$ is stored in receiver 1's cache. As we see in a moment, for $\alpha \in \{0, \frac{3D}{8}\}$, this allows to freely piggyback receiver 2's message $W_{d_2}^{(c1)}$ on receiver 1's message $W_{d_1}^{(u)}$ without compromising the rate to receiver 1.

Achievable rate-memory tuples: By Tuncel's seminal *broadcasting with side-information* result [13], communication in phase 1 (to both receivers) is very likely to be successful if the following two conditions hold:

$$R - \frac{2\mathcal{M}}{D} \leq F(1 - \delta_1)\beta_1 \quad (12a)$$

$$R \leq F(1 - \delta_2)\beta_1; \quad (12b)$$

communication in phase 2 is very likely to be successful if

$$R - \frac{2\mathcal{M}}{D} \leq F(1 - \delta_2)(1 - \beta_1). \quad (12c)$$

Inequalities (12) prove achievability of all rate-memory tuples $(R, \dots, R, \mathcal{M}, \dots, \mathcal{M})$, with R and \mathcal{M} satisfying (5) and

$$R \leq \begin{cases} \frac{4}{5}F(1 - \delta_1) + 2\frac{\mathcal{M}}{D}, & \text{if } \frac{\mathcal{M}}{R} \in [0, \frac{3D}{8}] \\ 2F(1 - \delta_1) + \frac{\mathcal{M}}{D} & \text{if } \frac{\mathcal{M}}{R} \in (\frac{3D}{8}, \frac{D}{2}]. \end{cases} \quad (13)$$

D. Discussion

Comparing the rate-memory tradeoffs in (8), (11) and (13), we see that it is advantageous to use unequal cache sizes and joint cache-channel coding. In particular, allowing larger caches at the weaker receivers (with higher packet erasure probabilities) both reduces the delivery-phase rates to the weaker receivers and increases rates to the stronger receivers.

IV. A JOINT CACHE-CHANNEL CODE FOR ARBITRARY DEMANDS

We now describe a joint cache-channel code that can be applied for any set of feasible demands \mathcal{D} , but we restrict attention to equal message rates

$$R_d = R, \quad d \in \{1, \dots, D\}.$$

We first treat the case where the K_0 weakest receivers (receivers 1 to K_0) have equal cache sizes and the remaining receivers do not have caches:

$$\mathcal{M}_k = \begin{cases} \mathcal{M} & \text{if } k \leq K_0 \\ 0 & \text{if } k > K_0 \end{cases}. \quad (14)$$

We explain later in Section IV-B how the scheme can be generalised to setups with unequal cache sizes.

A. Scheme for cache sizes satisfying (14)

Preliminaries: Choose a positive integer $t < K_0$, and let

$$\tau := \binom{K_0}{t}.$$

Split each message W_d into $(\tau + 1)$ -sub-messages,

$$W_d = (W_d^{(1)}, \dots, W_d^{(\tau+1)}),$$

where

$$W_d^{(i)} \in \{0, 1, \dots, 2^{nR(i)} - 1\}$$

and

$$R^{(i)} := \begin{cases} \frac{\mathcal{M}}{D \binom{K_0-1}{t-1}}, & \text{for } i = 1, 2, \dots, \tau \\ R - \frac{\mathcal{M}K_0}{Dt}, & \text{for } i = \tau + 1. \end{cases}$$

Caching Phase: Consider the K_0 weakest receivers. Let

$$\mathcal{R}_1, \mathcal{R}_2, \dots, \mathcal{R}_\tau$$

denote the τ different subsets of $\{1, \dots, K_0\}$ with size t . For each $i = 1, 2, \dots, \tau$, take the tuple

$$(W_1^{(i)}, W_2^{(i)}, \dots, W_D^{(i)})$$

and store it in the cache of each and every receiver in \mathcal{R}_i . Here we have stored $D \binom{K_0-1}{t-1}$ sub-messages in receiver k 's cache (for each $k \in \{1, 2, \dots, K_0\}$) with a total memory requirement

$$\left(2^{n \frac{\mathcal{M}}{D} \binom{K_0-1}{t-1}^{-1}}\right)^{D \binom{K_0-1}{t-1}} = 2^{n\mathcal{M}}.$$

Delivery phase: The demand tuple $\mathbf{d} \in \mathcal{D}$ is given, and we are required to communicate message W_{d_1} to receiver 1, W_{d_2} to receiver 2, and so on.

We consider sets of $(t+1)$ -receivers in $\{1, \dots, K_0\}$. Within these sets, each subset of t receivers shares a sub-message that is demanded (but unknown) by the remaining $(t+1)$ -th receiver. For each set of $(t+1)$ -receivers, we form the ‘‘XOR’’ of the $(t+1)$ sub-messages having the two above mentioned properties, that is, being known at t of the receivers and demanded by the remaining $(t+1)$ -th receiver. For example,

for the subset of receivers $\{1, \dots, t+1\}$, we form the XOR message

$$\bigoplus_{k=1}^{t+1} W_{d_k}^{(i_k)},$$

where the addition is modulo $2^{\mathcal{M}/(D \binom{K_0-1}{t-1})}$ (or, equivalently, a bitwise XOR operation); and for each $k \in \{1, \dots, t+1\}$, i_k is such that

$$\mathcal{R}_{i_k} \triangleq \{1, \dots, k-1, k+1, \dots, t+1\}. \quad (15)$$

Notice that (15) implies that $W_{d_k}^{(i_k)}$ is stored in the caches of receivers $1, \dots, k-1, k+1, \dots, t$, but not at receiver k .

We use a time-sharing scheme to send the XOR messages as well as all other messages to be transmitted. The time-sharing comprises K phases. Each phase $k \in \{1, \dots, K\}$ is constructed so that it can be decoded by Receivers $k, k+1, \dots, K$. Phase k is of length $\beta_k n$, where

$$\sum_{k=1}^K \beta_k = 1, \quad 0 \leq \beta_k \leq 1. \quad (16)$$

In phase $k \in \{1, \dots, K_0\}$, we send

- the XOR messages that are demanded by receiver k but not by receivers 1 to $k-1$;
- the uncached message $W_{d_k}^{(\tau+1)}$ demanded by Receiver k ;
- the first $nC_{k,\tilde{k}}$ bits of sub-messages $W_{d_{\tilde{k}}}^{(i)}$, for every $\tilde{k} \in \{K_0+1, \dots, K\}$ and every $i \in \{1, \dots, \tau\}$ such that $k \in \mathcal{R}_i$. These messages are all known to receiver k and therefore do not limit the decoding at receiver k . The rates $\{C_{k,\tilde{k}}\}$ are parameters of a scheme. As we shall see, when they are chosen sufficiently small, but positive, and $\delta_{k+1} < \delta_k$, then sending these bits does not limit the decoding at receivers $k+1, k+2, \dots, K$. In fact, similarly to our motivating example, in this case, the transmitted bits of sub-messages $W_{d_{\tilde{k}}}^{(i)}$ can be freely piggybacked on the other messages transmitted in this phase k .

In phase $k \in \{K_0+1, \dots, K\}$ we send:

- the sub-messages of W_{d_k} that have not been sent in any previous phase.

Achievable rate-memory tuples:

Proposition 2: A rate-memory tuple $(R, \dots, R, \mathcal{M}_1 = \mathcal{M}, \dots, \mathcal{M}_{K_0} = \mathcal{M}, 0, \dots, 0)$ is achievable if for some

- positive integer t ;
- nonnegative K -tuple $(\beta_1, \dots, \beta_K)$ satisfying (16); and
- nonnegative real numbers $\{C_{k,\tilde{k}}\}$ with $k \in \{1, \dots, K_0\}$ and $\tilde{k} \in \{K_0+1, \dots, K\}$;

the following $(K+K_0)$ -conditions in (17) hold.

- 1) For each $k \in \{1, \dots, K_0 - t - 1\}$, we have

$$R \leq F(1 - \delta_k) + \frac{\mathcal{M}}{D \binom{K_0-1}{t-1}} \left(\binom{K_0}{t} - \binom{K_0-k}{t} \right) \quad (17a)$$

and

$$R + \sum_{\tilde{k}=K_0+1}^K C_{k,\tilde{k}} \leq F(1 - \delta_{k+1})$$

$$+ \frac{\mathcal{M}}{D \binom{K_0-1}{t-1}} \left(\binom{K_0}{t} - \binom{K_0-k}{t} \right). \quad (17b)$$

2) For each $k \in \{K_0 - t, \dots, K_0\}$, we have

$$R \leq F(1 - \delta_k) + \frac{\mathcal{M}K_0}{Dt} \quad (17c)$$

and

$$R + \sum_{\tilde{k}=K_0+1}^K C_{k,\tilde{k}} \leq F(1 - \delta_{k+1}) + \frac{\mathcal{M}K_0}{Dt}. \quad (17d)$$

3) Finally, for each $k \in \{K_0 + 1, \dots, K\}$, we have

$$R - \sum_{k'=1}^{K_0} C_{k',k} \leq F(1 - \delta_k). \quad (17e)$$

Proof outline: For each $k \in \{1, \dots, K_0 - t\}$, Condition (17a) ensures that receiver k can reliably decode the sub-messages sent during phase k , and Condition (17b) ensures that all of the stronger receivers in $\{k + 1, \dots, K\}$ can also reliably decode these sub-messages. Similarly, Condition (17c) ensures that each receiver $k \in \{K_0 - t, \dots, K_0\}$ can reliably decode the sub-messages sent in phase k , and Condition (17d) ensures that all of the stronger receivers in $\{k + 1, \dots, K\}$ can also reliably decode these sub-messages. Finally, Condition (17e) ensures that each receiver $k \in \{K_0 + 1, \dots, K\}$ can decode the sub-messages sent in phase k . ■

Discussion: The parameters $\{C_{k,\tilde{k}}\}$ describe the gain that our scheme achieves over separate cache-channel coding schemes. If some of these rates are strictly larger than 0, then our scheme strictly outperforms separate cache-channel coding. It is possible to choose them strictly positive whenever the erasure probabilities $\delta_1, \dots, \delta_{K_0}$ are not all equal.

We took advantage of the fact that receiver k has already cached the additional $nC_{k,\tilde{k}}$ message bits that are sent in phase k . Some of these bits are also available to the next-stronger receivers $k + 1, k + 2, \dots$. For simplicity, we ignored this fact in our analysis, and it is likely that further gains can still be made.

B. Scheme for unequal cache sizes

Assume now that

$$\mathcal{M}_1 \geq \mathcal{M}_2 \geq \dots \mathcal{M}_K \geq 0. \quad (18)$$

Our scheme in the previous subsection is easily extended to this more general setup using time-sharing. Specifically: Let β_1, \dots, β_K be real numbers in the interval $[0, 1]$ that sum up to 1. Over a fraction of time β_i , $i \in \{1, \dots, K\}$, we use our scheme in the previous subsection assuming that only the first $K_0^{(i)} = K + 1 - i$ receivers have caches of equal cache size $\mathcal{M}^{(i)} = \beta_i^{-1}(\mathcal{M}_{K-i+1} - \mathcal{M}_{K-i+2})$. (Set $\mathcal{M}_{K+1} \triangleq 0$.)

V. SINGLE COMMON DEMAND

In this section we consider the optimistic case where all receivers demand the same message. This corresponds to

$$\mathcal{D} = \{(d_1, \dots, d_K) \in \{1, \dots, D\}^K : d_1 = d_2 = \dots = d_K\}.$$

The rates R_1, \dots, R_D can be arbitrary, i.e., do not have to be equal as in the previous section.

A. Result

Theorem 3: A rate-memory tuple $(R_1, \dots, R_D, \mathcal{M}_1, \dots, \mathcal{M}_K)$ is achievable if and only if,

$$R_d \leq \min_{k \in \{1, \dots, K\}} ((1 - \delta_k)F + \mathcal{M}_{k,d}), \quad d \in \{1, \dots, D\}, \quad (19)$$

for some nonnegative numbers $\{\mathcal{M}_{k,d}\}$ that satisfy

$$\sum_{d=1}^D \mathcal{M}_{k,d} \leq \mathcal{M}_k, \quad k \in \{1, \dots, K\}. \quad (20)$$

Proof: See the following two subsections. ■

We thus again wish to allocate small cache sizes to strong receivers and large cache sizes to weak receivers.

If we used separate cache-channel codes, Constraint (19) is replaced by

$$\max_{k \in \{1, \dots, K\}} (R_d - \mathcal{M}_{k,d}) \leq \min_{k \in \{1, \dots, K\}} (1 - \delta_k)F, \quad (21)$$

and the benefit of having unequal cache sizes $\{\mathcal{M}_d\}$ at the different receivers disappears.

B. Proof of achievability

We propose the following scheme.

Caching phase: Each receiver k stores in its cache the first $n\mathcal{M}_{k,d}$ bits of each Message W_d , for $d \in \{1, \dots, L\}$, where

$$\sum_{d=1}^N \mathcal{M}_{k,d} \leq \mathcal{M}_k, \quad (22)$$

in order to satisfy the cache-memory constraint.

Delivery phase: Assume $d_1 = d_2 = \dots = d_K = d^*$. Use an i.i.d. Bernoulli-1/2 point-to-point code to send Message W_{d^*} to all receivers. Each receiver k knows the first $n\mathcal{M}_{k,d^*}$ bits of this message, and thus during its decoding it can restrict attention to the part of the codebook that corresponds to these bits. For receiver k it is thus as if the transmitter had sent only its missing bits over the channel.

Alternatively, a joint cache-channel code based on Tuncel's virtual binning technique [13] can be used for the delivery phase.

Achievable rate-memory tuples: By [13], whenever

$$R_d - \mathcal{M}_{k,d} \leq (1 - \delta_k)F, \quad \forall d \in \{1, \dots, D\}, \quad k \in \{1, \dots, K\}, \quad (23)$$

the probability of error can be made arbitrarily small as $n \rightarrow \infty$.

C. Proof of Converse

Fix a block length n , and define

$$\mathcal{M}_{k,d} \triangleq \frac{1}{n} I(W_d; \mathbb{Z}_k), \quad k \in \{1, \dots, K\}, \quad d \in \{1, \dots, D\}. \quad (24)$$

For each $k \in \{1, \dots, K\}$,

$$\sum_{d=1}^N \mathcal{M}_{k,d} = \sum_{d=1}^D \frac{1}{n} I(W_d; \mathbb{Z}_k)$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{d=1}^D (H(W_d) - H(W_d|\mathbb{Z}_k)) \\
&= \frac{1}{n} (H(W_1, \dots, W_D) - \sum_{d=1}^D H(W_d|\mathbb{Z}_k)) \\
&\leq \frac{1}{n} (H(W_1, \dots, W_D) - H(W_1, \dots, W_N|\mathbb{Z}_k)) \\
&= \frac{1}{n} I(W_1, \dots, W_D; \mathbb{Z}_k) \\
&\leq \frac{1}{n} H(\mathbb{Z}_k) \leq \mathcal{M}_k,
\end{aligned} \tag{25}$$

where the second and fourth equalities follow by the definition of mutual information; the third equality because the messages are independent; the first inequality because the sum of marginal entropies of a tuple of random variables, is at least as large as the joint entropy of this tuple.

In the following, let ϵ_n denote any sequence that tends to 0 as $n \rightarrow \infty$. Fix an achievable rate-memory tuple $(R_1, \dots, R_N, \mathcal{M}_1, \dots, \mathcal{M}_K)$. Also for an arbitrary large n , let $\{\mathbb{Z}_1, \dots, \mathbb{Z}_K\}$, $\{f_d\}$, and $\{\varphi_{k,d}\}$ denote cache content, encoding functions, and decoding functions achieving this rate-memory tuple. Fix now $d^* \in \{1, \dots, D\}$ and $k \in \{1, \dots, K\}$, and let $X^n = f_{d^*}(W_1, \dots, W_D)$ and Y^n denote inputs and outputs corresponding to demand $\mathbf{d}^* \triangleq (d^*, d^*, \dots, d^*)$. We have

$$\begin{aligned}
R_{d^*} &\leq \frac{1}{n} H(W_{d^*}) \\
&= \frac{1}{n} I(W_{d^*}; Y_k^n, \mathbb{Z}_k) + \frac{1}{n} H(W_{d^*} | Y_k^n, \mathbb{Z}_k) \\
&\leq \frac{1}{n} I(W_{d^*}; Y_k^n | \mathbb{Z}_k) + \frac{1}{n} I(W_{d^*}; \mathbb{Z}_k) + \epsilon_n \\
&= \frac{1}{n} \sum_{t=1}^n (H(Y_{k,t} | \mathbb{Z}_k, Y_k^{t-1}) - H(Y_{k,t} | W_{d^*}, Y_k^{t-1}, \mathbb{Z}_k)) \\
&\quad + \mathcal{M}_{k,d^*} + \epsilon_n \\
&\leq \frac{1}{n} \sum_{t=1}^n (H(Y_{k,t}) - H(Y_{k,t} | X_{k,t})) + \mathcal{M}_{k,d^*} + \epsilon_n \\
&= \frac{1}{n} \sum_{t=1}^n I(Y_{k,t}; X_{k,t}) + \mathcal{M}_{k,d^*} + \epsilon_n \\
&\leq (1 - \delta_k) F + \mathcal{M}_{k,d^*} + \epsilon_n,
\end{aligned} \tag{26}$$

where the second inequality follows by Fano's inequality; the third inequality because conditioning cannot increase entropy and because of the Markov chain $(W_{d^*}, Y_k^{t-1} | \mathbb{Z}_k) \rightarrow X_t \rightarrow Y_{k,t}$; and the last inequality by the capacity of the erasure channel; all equalities follow by the definition and the chain rule of mutual information.

Letting $n \rightarrow \infty$, and thus $\epsilon_n \rightarrow 0$, establishes the desired converse.

REFERENCES

- [1] M. A. Maddah-Ali, U. Niesen, "Fundamental limits of caching," in *IEEE Trans. on Inf. Theory*, vol. 60, no. 5, pp. 2856–2867, May 2014.
- [2] M. A. Maddah-Ali, U. Niesen, "Decentralized coded caching attains order-optimal memory-rate tradeoff," *IEEE/ACM Trans. on Networking*, vol. PP, no. 1, pp. 1, 2014.
- [3] U. Niesen and M. A. Maddah-Ali, "Coded caching with nonuniform demands," in *Proc. IEEE INFOCOM Workshop*, 2014.
- [4] S. Wang, X. Tian and H. Liu, "Exploiting the unexploited of coded caching for wireless content distribution," in *Proc. IEEE ICNC*, 2015.
- [5] R. Pedarsani, M. A. Maddah-Ali and U. Niesen, "Online coded caching," *IEEE/ACM Trans. on Networking*, vol. PP, no.1, pp. 1, 2015.
- [6] N. Karamchandani, U. Niesen, M. A. Maddah-Ali, S. Diggavi, "Hierarchical coded caching," *submitted to IEEE Trans. on Inf. Theory*. Online: <http://arxiv.org/pdf/1403.7007>.
- [7] J. Hachem, N. Karamchandani, and S. Diggavi, "Content caching and delivery over heterogeneous wireless networks." Online: <http://arxiv.org/pdf/1404.6560>.
- [8] H. Ghasemi and A. Ramamoorthy, "Improved lower bounds for coded caching," Online: <http://arxiv.org/pdf/1501.06003>.
- [9] C.-Y. Wang, S. H. Lim, and M. Gastpar, "Information-theoretic caching: sequential coding for computing," *submitted to the IEEE Trans. on Inf. Theory*. Online: <http://arxiv.org/abs/1504.00553v1>.
- [10] M. Ji, A. M. Tulino, J. Llorca, and G. Caire, "Order-optimal rate of caching and coded multicasting with random demands," *submitted to the IEEE Trans. on Inf. Theory*. Online: <http://arxiv.org/abs/1502.03124>.
- [11] W. Huang, S. Wang, L. Ding, F. Yang, and W. Zhang, "The performance analysis of coded cache in wireless fading channel," *submitted to GC2015*. Online: <http://arxiv.org/abs/1504.01452>.
- [12] R. Urbanke and A. Wyner "Packetizing for the erasure broadcast channel with an internet application," 1999. Online: <http://lthcwww.epfl.ch/~ruediger/papers/inft.ps>.
- [13] E. Tuncel, "Slepian-Wolf coding over broadcast channels," *IEEE Trans. on Inf. Theory*, vol. 52, no. 4, pp. 1469–1482, 2006.